

《BERT 的逻辑化》

3.0 修正版

YKY

October 17, 2020

1 Summary

- BERT 在训练过程中「被逼」预测遮掩的词语，由此诱导出人类的知识，这已经具有 AGI 的雏形
- 从经典逻辑 AI 的角度看，BERT 内部可能有符合逻辑思维的结构，而我们很惊讶地发现，透过 Curry-Howard 对应，BERT 可以看成是一种另类的逻辑
- AGI 系统可以和 逻辑结构（以 范畴论、topos 表述）建立紧密的联系，这个联系可以指导往后的 AGI 发展路线，非常方便
- 我需要一些合作者帮助，特别是 BERT 和 soft Actor-Critic 方面

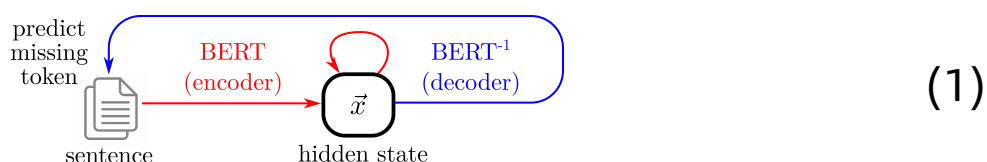
Contents

2	BERT's ground-breaking significance: closed-loop training	4
3	BERT's internal architecture	4
4	Symmetry in logic	5
5	Symmetric neural network	6
6	Knowledge graphs	6
7	Logicalization of BERT	7
8	Connection between AI and logic	7
9	What is attention?	8
10	Predicates vs propositions	8
11	"Attention is all you need" ?	9
12	Why is BERT so successful?	10
13	Modifying BERT with an attention-like mechanism	11

14 Abductive interpretation of natural language	11
15 Capturing semantics more broadly	12
16 Content-addressable long-term memory	12
17 Doubts about logicism	13
18 References	13

2 BERT 的革命性意义：「闭环路训练」

- BERT 利用平常的文本 induce 出知识，而这 representation 具有 通用性 (universality)：



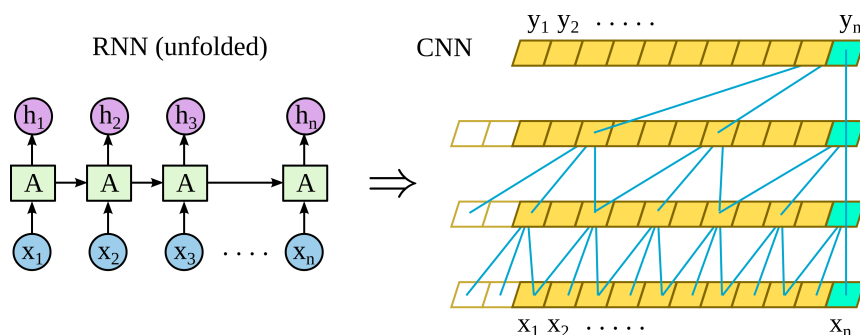
换句话说：隐状态的 representation 压缩了句子的意思，而它可以应用在别的场景下

- This implies that human-level AI can be *induced* from existing corpora, 而不需重复 像人类婴儿成长的学习阶段
- 这种训练方法是较早的另一篇论文提出，它并不属于 BERT 的内部结构

3 BERT 的内部结构

其实，BERT 也是混合了很多技巧 发展而成的：

- BERT 基本上是一个 seq-to-seq 的运算过程
- Seq-to-seq 问题最初是用 RNN 解决的
- 但 RNN 速度较慢，有人提出用 CNN 取代：



- CNN 加上 attention mechanism 变成 Transformer
- 我的想法是重复这个思路，但引入 逻辑的对称性

4 Symmetry in logic

- 词语 组成 句子，类比於 逻辑中，概念 组成 逻辑命题
- 抽象地说，逻辑语言 可以看成是一种有 2 个运算的 代数结构：加法 (\wedge , 合并命题, 可交换) 和 乘法 (\cdot , 用作概念合成, 不可交换)
- 例如：

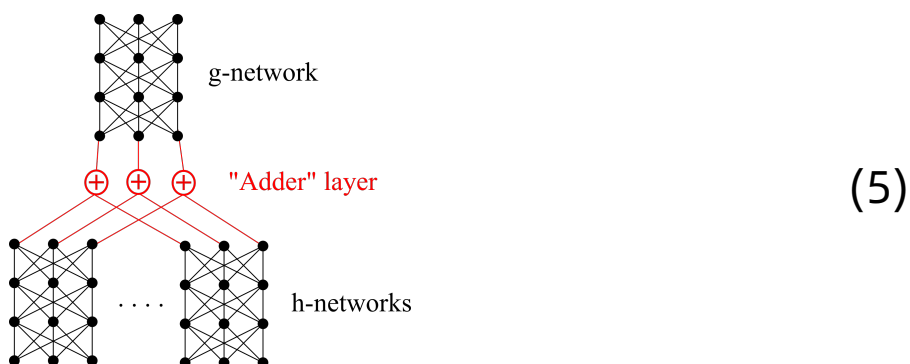
$$\begin{aligned} A \wedge B &\equiv B \wedge A \\ \text{下雨} \wedge \text{失恋} &\equiv \text{失恋} \wedge \text{下雨} \end{aligned} \quad (3)$$

- Word2Vec 也是革命性的；由 Word2Vec 演变成 Sentence2Vec 则比较容易，基本上只是 向量的 合并 (concatenation); Sentence 对应於 逻辑命题
- 但 命题的 集合 需要用 symmetric NN 处理，因为 集合的元素 是 顺序 无关的

5 Symmetric neural network

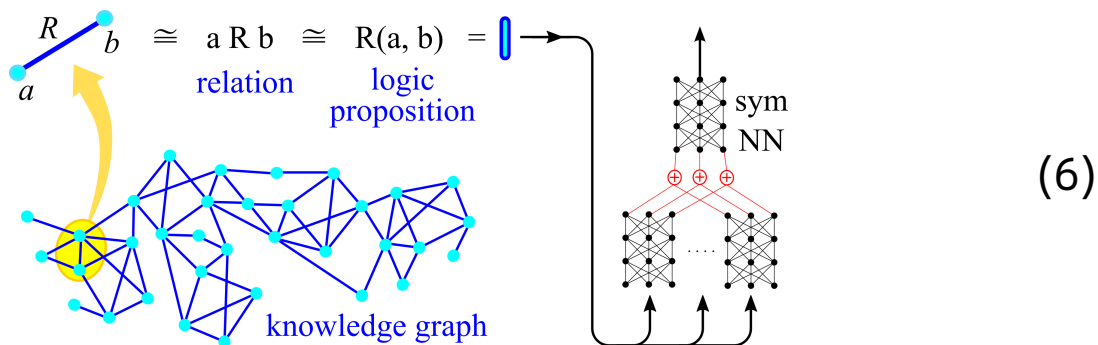
- Symmetric NN 问题 已经由 两篇论文解决了:
[PointNet 2017] [DeepSets 2017]
- Any symmetric function can be represented by the following form
(a special case of the Kolmogorov-Arnold representation of functions):

$$f(x, y, \dots) = g(h(x) + h(y) + \dots) \quad (4)$$



6 知识图谱 (knowledge graphs)

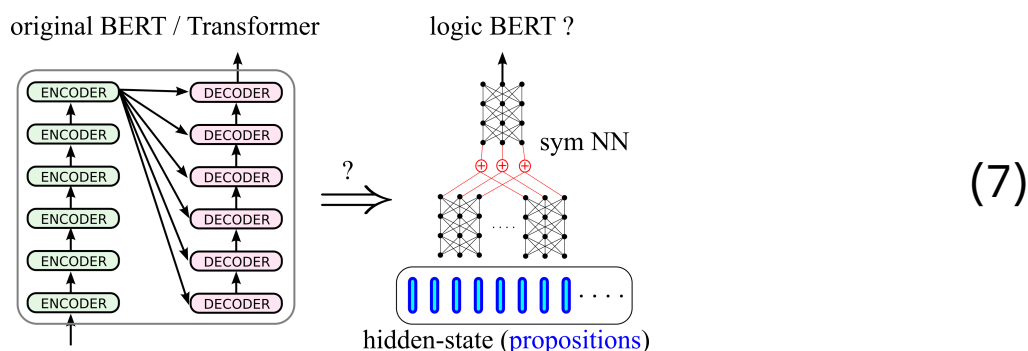
- 知识图谱 不能直接输入神经网络,它必需分拆成很多 edges,每个 edge 是一个 关系, 也是一个 逻辑命题; 也可以说 "graphs are isomorphic to logic"



- 而这些 edges 似乎必需用 symmetric NN 处理, 因为它们是 permutation invariant
- 逻辑化 建立了 知识图谱 和 BERT 之间的一道桥梁
接下来讨论 BERT....

7 BERT 的逻辑化

- 可以强逼 BERT 的 隐状态 变成 “set of propositions” 的形式，方法是将 对称性 施加在 Encoder 上：



- 下面会看到，BERT 的「注意力机制」已经是对称的，它可以做 逻辑推导

8 逻辑与 AI 之间的联系

- 既然 AI 基於 逻辑，则 AI 与逻辑之间 必然存在 精确 (precise) 的联系
- BERT 似乎是在执行 句子之间的变换，而这些句子是 word embedding 的 concatenation，例如：

$$\text{苏格拉底} \cdot \text{是} \cdot \text{人} \xrightarrow{BERT} \text{苏格拉底} \cdot \text{会} \cdot \text{死} \quad (8)$$

这个做法看似很「粗暴」，其实它和 逻辑式子 的作用一样：

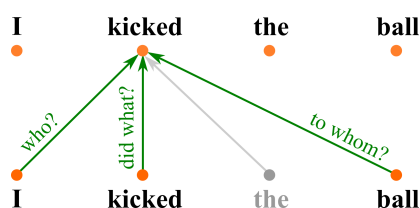
$$\forall x. \text{Human}(x) \Rightarrow \text{Mortal}(x) \quad (9)$$

而这式子，根据 Curry-Howard 对应，就是 (8) 的函数映射！

- 我会在另一辑 slides 里 简介一下这些理论；可以说，逻辑的 几何结构 是「永恒」的，它可以指示 AI 的 长远发展

9 Attention 是什么?

- 注意力 最初起源於 Seq2seq, 后来 BERT 引入 self-attention
- Attention 的本质就是 加权, 权值 可以反映 模型 关注的点
- For each input, attention weighs the **relevance** of every other input and draws information from them accordingly to produce the output
- 在 BERT 里, attention 是一种 **words** 之间的关系:

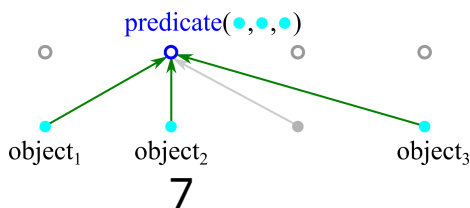


(10)

- 但, 从逻辑的角度看, word \neq 命题
- 在逻辑学上, 必需分清 命题**内部** 与 命题**之间** 这两个层次, 非常关键!

10 谓词 (predicates) vs 命题 (propositions)

- “Predicate” 来自拉丁文「断言」的意思
- 逻辑里, predicate 代表一个 没有主体 / 客体的断言, 换句话说, 是一个有「洞」的命题
- **命题** = **谓词** (predicate) + **主体 / 客体** (统称 objects)
- 例如: Human(John), Loves(John, Mary)
- 从逻辑的角度看, attention 的输出可以看成是 predicate 和 objects 的**结合**:



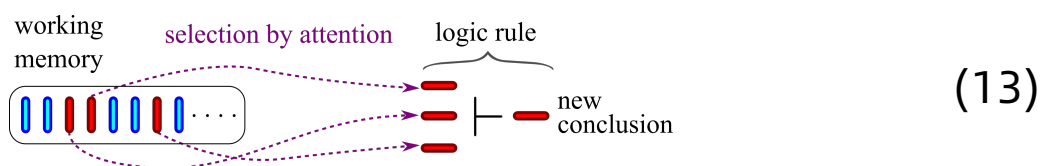
(11)

- 形象地说：

$$\begin{array}{ccc} \text{predicate} + \text{objects} & = & \text{proposition} \\ \circ + \bullet \bullet \bullet \dots & = & \text{—} \end{array} \quad (12)$$

11 “Attention is all you need” ?

- 类似地，**高层的** attention 可以处理 **命题之间** 的关系
- 我们希望 attention 做到的是 **选择** 有关联的命题，去做逻辑推导：



- 但从 N 个命题中选择 K 个，可以有 $\binom{N}{K}$ 个子集，是 exponential 的
- BERT 的做法是：每次只输出 N 个命题，而前提的 “support” 也是上一层的 全部 N 个命题，每个前提的「影响力」由 matrix 权重决定
- 根据 Curry-Howard isomorphism, BERT 的映射 其实对应於某种 **另类的逻辑** (BERT 的设计者可能也意识到这点)，这种逻辑的好处是运行**非常快**
- BERT 的逻辑 看似有局限，但这种表面的局限未必阻止它是 universal 的逻辑
- 关键是在 速度 与 逻辑的 expressive power 之间 找到平衡
- 最简单的 attention 公式是 (其中 Q, K, V = query, key, value 矩阵):

$$\mathbf{y}_j = \sum_i \langle Q \mathbf{x}_j, K \mathbf{x}_i \rangle V \mathbf{x}_i \quad (14)$$

(**红色** 代表注意力的 focus) 这对应於一个逻辑式子：

$$\equiv x_1 \wedge x_2 \wedge x_3 \wedge x_4 \vdash y_2 \quad (15)$$

亦即是说： y_j 是由 x_1, \dots, x_n 得出的逻辑结论 with focus on x_j

- 所谓 “focus” 并不是 逻辑概念，它只是 BERT 加速的 heuristic
- 容易看到，attention 对於 x_1, \dots, x_n 是 交换不变的 (equivariant)，这表示 每层 attention 的输出是一些 逻辑命题，和我的理论相符
- Multi-head attention 亦有一个很好的逻辑解释：即使 focus = x_j ，亦有其他不同的前提，可以导致不同的结论，例如：

$$x_1 \wedge x_2 \wedge x_3 \vdash y_1 \quad , \quad x_1 \wedge x_4 \wedge x_5 \vdash y_2 \quad (16)$$

12 BERT 为什么成功？

- 6 层的 BERT 有 $512 \times 6 = 3072$ 个 head，如果 $8 \times$ multi-head 则有 24576 个
- Each head does not simply correspond to 1 formula in conventional logic, may require further in-depth analysis....
- 我的猜测是：BERT 的高层 representation 是一些有 “high-level” 意义的命题，就像在视觉中，高层特征 代表一些复杂的物体
- The embedding of high-level propositions in vector space may be “semantically dense”, meaning that slight changes in the vector position may convey many different meanings
- 由於 logic rules 是以 6 层的 hierarchy 组织而成，这结构具有 “deep learning” 的特性

13 改良 BERT: 类似 attention 的方法

- The BERT attention formula (14) has some unnecessary restrictions, where generally we just need a symmetric function in the x_i 's
- The general form of symmetric functions is given by (4)
- Immitating BERT, we introduce a “focus” of attention on x_j :

$$y_j = g(h(\mathbf{x}_j, \mathbf{x}_1) + \dots + h(\mathbf{x}_j, \mathbf{x}_n)) \quad (17)$$

this preserves **equivariance**

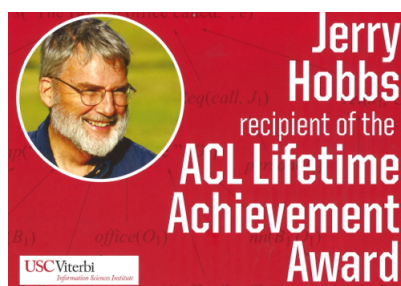
- We can use this function to replace the entire BERT Encoder:



14 逆因推理 (abduction) 与 自然语言理解

- 根据 Jerry Hobbs 的 “abductive interpretation of natural language” 理论，语言理解 是一个 **解释** (explain) 的过程
- **解释** 在逻辑上等同於 **逆因推理** (abduction), 亦即是 **逻辑蕴涵** (implication, $A \Rightarrow B$) 的反方向
- 举例来说：天气热 \Rightarrow 流汗，所以「天气热」就是「流汗」的 **解释**
- 这个理论 今天很少人知道，因为属於 经典逻辑 AI 时期

2013 年 他
获得
计算语言学
终身成就奖



15 捕捉更广泛的 semantics

- 阅读时，我们（人脑）有时可以预测下个 word（这是 BERT 训练的目标），但有时即使不能预测，但看到 next word 之后，仍会有某种「意料之内」的感觉

- 举例来说：

「天气热，我不停 流汗」

「天气热，我不停 吃冰淇淋」

第二个例子是比较少见的，但也合理

- （从逻辑角度看）BERT 只会预测 **最有可能** 的结论：

$$\text{前提} \xrightarrow{BERT} \text{预测} \quad (19)$$

其实我们想要的是：

$$\text{前提} \longrightarrow \text{很多不同的 预测} \quad (20)$$

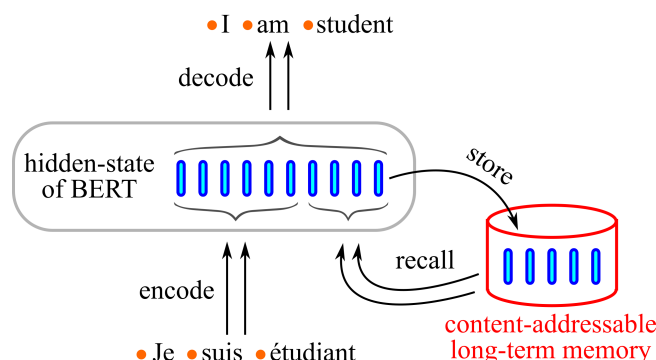
而只要其中一个预测成功，即给予「奖励」

- 这个情况和 **强化学习** 中的 “stochastic actions” 类似，需要的是：stochastic, multi-modal, continuous actions（例如 SAC, soft actor-critic）

16 Content-addressable long-term memory

- Content-addressable memory 的想法来自 Alex Graves *et al* 的 Neural Turing Machine [2014]

- 以前 BERT 的隐状态 没有逻辑结构，我们不是很清楚它的内容是什么；逻辑化之后，BERT 内部的命题可以储存在 长期记忆 中：



(21)

- 这种系统 已非常接近 strong AI，而这是有赖 逻辑化 才能做到的
- 但这个 idea 暂时仍未成熟，有待更多研究

17 对 逻辑主义 的质疑

- 很多人怀疑：人脑真的用 逻辑 思考吗？
- 其实我们每句表达的 语言，都是逻辑形式的 (logical form)
- 直觉认为，人脑 构造一些 models，再从 model 中「读出」一些结论
- 例如给定一个描述：「已婚妇人出轨，用刀刺死丈夫」



(22)

- 那么 妻子穿著什么衣服？衣服什么颜色？这些都是 臆想 出来的细节，是不正确的
- 这个 model 可以有哪些细节？答案是：任何细节都不可以有，除非是逻辑上蕴含的，或被 逻辑约束
- Model 本身可以是一些 抽象的逻辑命题 构成的，这也合理；反而，一个有很多感官细节的 model 并不合理
- 其实人脑可能 比我们想像中 更接近逻辑

18 References

欢迎提问和讨论 😊